# The Study of Intelligent Algorithm in Particle Identification of Heavy-Ion Collisions at Low and Intermediate Energies[*]

Gao-Yi Cheng,[1, 2, 3] Qian-Min Su,[1, †] Xi-Guang Cao,[2, 3, ‡] and Guo-Qiang Zhang[2, 3]

[1]*School of Electronic and Electrical Engineering, Shanghai University of Engineering Science, Shanghai 201620, China*

[2]*Shanghai Advanced Research Institute, Chinese Academy of Sciences, Shanghai 201210, China*

[3]*Shanghai Institute of Applied Physics, Chinese Academy of Sciences, Shanghai 201800, China*

Traditional particle identification methods are time consuming, experience-dependent, and poor repeatability challenges in heavy-ion collisions at low and intermediate energies. Researchers urgently need solutions to the dilemma of traditional particle identification methods. This study explores the possibility of applying intelligent learning algorithms to the particle identification of heavy-ion collisions at low and intermediate energies. Multiple intelligence algorithms, including XgBoost and TabNet, were selected to test datasets from the neutron ion multi-detector for reaction-oriented dynamics (NIMROD–ISiS) and Geant4 simulation. Machine learning algorithms based on tree structures and deep learning algorithms e.g. TabNet show excellent performance and generalization ability. Adding additional data features besides energy deposition can improve the algorithm's identification ability when the data distribution is nonuniform. Intelligent learning algorithms can be applied to solve the particle identification problem in heavy-ion collisions at low and intermediate energies.

Keywords: Heavy-ion collisions at low and intermediate energies, Machine learning, Ensemble learning algorithm, Particle identification, Data imbalance

## I. INTRODUCTION

Intelligent algorithms play crucial roles in nuclear physics. Challenges in nuclear physics experiments include high complexity, extensive data, time-consuming experiments, and intricate models. Taking particle collision experiments as an example, millions of terabytes of data are generated daily for heavy-ion collisions at high energies. Therefore, the extraction of useful information from complex experimental data has become an enormous challenge.

Large-scale experiments such as ATLAS, ALICE, and CMS have already applied machine-learning and deep-learning algorithms [1–4] to analyze and process experimental data. Typical examples include research on the particle–track reconstruction problem [5–8] in high-energy physics experiments, data analysis, and pattern recognition of the Higgs boson [9–13]. The application of machine learning in particle physics can be seen in a large-scale dynamic review [https://iml-wg.github.io/HEPML-LivingReview/] and the website opened by the ML Physics Portal [14–17].

Currently, research on intelligence algorithms in nuclear physics experiments [18–21] focuses on data analysis, such as the masses of atomic nuclei [22–26], nuclear charge radii [27–31], decay half-lives [32–37], critical reaction thresholds [38], and spallation reaction cross-sections [39], etc. In addition to using machine learning algorithms to investigate various physical issues [40–42], researchers have used these algorithms to analyze experimental data [43–45]. This involves tasks, such as

[†] Corresponding author, Qian-Min Su, suqm@sues.edu.cn

[‡] Corresponding author, Xi-Guang Cao, caoxg@sari.ac.cn

particle trajectory reconstruction, vertex reconstruction [46], and particle identification in nuclear reactions. Advancements in experimental equipment and related technologies have facilitated the integration of machine learning and nuclear physics.

Current research on particle identification focuses on high-energy particle physics. To date, research on particle identification has mainly focused on identifying particle types [47] and separating rare particles from background signals. The data and algorithms used for particle identification depend on the type of detector. For example, the output data of a calorimeter detector can be processed and converted into matrix data; therefore, image algorithms, such as CNN and GNN, can be used for processing. The research and applications of machine learning in particle identification have mainly focused on LHC detectors, such as calorimeters [48–50] and Cherenkov detectors [51]. Moreover, a new research focus in recent years on LHC experiments has been the development of new detector software and hardware based on machine-learning and deep-learning algorithms [52].

Compared with other nuclear reactions, the particles generated in heavy-ion collisions at low and intermediate energies are of various types and have complex energy distributions. Numerous fragments have similar charges and masses. Experiments on heavy-ion collisions depend on the energy resolution of the detector and require a detection array with large solid-angle coverage. Therefore, the identification of dozens or even hundreds of reaction products from independent detection units is challenging. Traditional particle identification methods include telescope [53], time-of-flight [54], magnetic spectrometer, Bragg spectroscopy, and pulse shape discrimination methods. These methods are often combined to improve identification ability, especially for heavy fragments with minor differences in charge and mass numbers between adjacent fragments. The performance of the traditional methods for heavier particles is hindered by their dependence on experience, poor repeatability, and time consumption. The precise identification of charge and mass numbers is fundamental to all research related to heavy-ion collisions, and is a very powerful method for studying exotic nuclear configurations [55–60]. Compared with particle identification in particle physics, the wide variety and slight differences in the charge and mass numbers of charged particles produced in heavy-ion reactions pose significant challenges for existing particle identification methods. Therefore, the development of a universal, efficient, and high-precision particle identification method based on machine learning techniques will significantly boost the study of heavy-ion collisions.

Parker et al. [61] devised a 5-layer neural network and evaluated its performance on the 22nd and 23rd detectors of a neutron ion multi-detector for reaction-oriented dynamics (NIMROD–ISiS). We also used a dataset from NIMROD–ISiS detector array. This study aimed to identify the particle charge and mass numbers in heavy-ion collisions at low and intermediate energies. Supervised learning algorithms were used to train particle identification models based on $\Delta$E-E energy deposits from telescope (or super-telescope) detectors in heavy-ion collisions. Machine learning and deep learning algorithms were applied to identify the particles' charge and mass numbers, and their capabilities were compared.

## II. DATASET AND METHODS

Real-world data (RWD) come from experiments on heavy-ion collisions at low and intermediate energies carried out at the Cyclotron Institute of Texas A&M University and consist of reaction products detected by the NIMROD–ISiS array [62, 63]. The NIMROD–ISiS detector array comprised 14 rings. Experimental data were obtained from 143 detectors, including 124 telescope detectors and 19 super telescope detectors with ring numbers ranging from 2 to 15. The detection system included Si detectors and CsI (Tl) scintillators with angles ranging from $3.6°$ to $167.0°$. The back half of the NIMROD ($90.0°$ - $167.0°$) consists of half the Indiana Silicon Sphere. Si detectors were combined with the CsI detectors as 'telescopes,' while some were equipped with two Si detectors in tandem, known as 'super telescopes' ($3.6°$-$45°$), enhancing the ability to identify mass

52 numbers of heavier fragments. The capacity to include ionization chambers in front of Si detectors is also available. Fig. 1

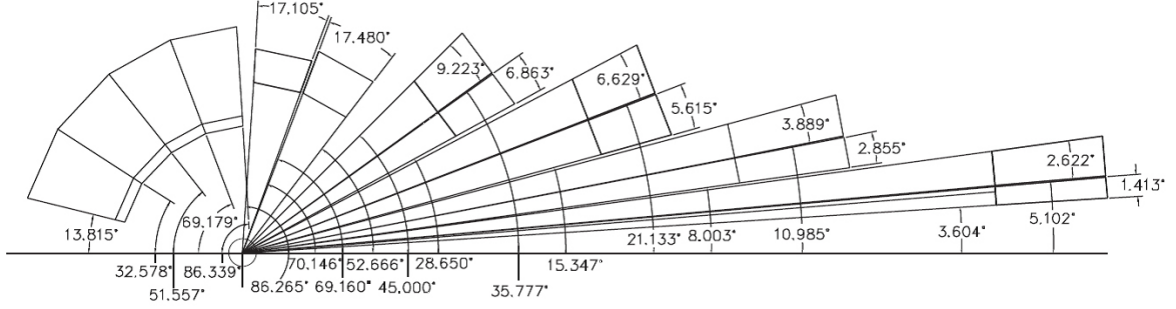53 shows the structure of the NIMROD-ISIS.



Fig. 1. Schematic diagram of the NIMROD–ISiS detector array layout, from Texas A&M NIMROD–ISiS official website (https://cyclotron.tamu.edu/nimrod/)

54 In addition to the dataset from Texas A&M University, Geant4 [64] was used to simulate heavy-ion collisions at intermediate

55 energies. QMD model with G4IonQMDPhysics was used as an event generator to simulate the reaction process of a beam inci-

56 dent on a target. The detection processes in Geant4 include electromagnetic interactions (G4EmStandardPhysics), energy trans-

57 fer and loss (G4EmExtraPhysics and G4StoppingPhysics), decay processes (G4DecayPhysics and G4RadioactiveDecayPhysics),

58 and elastic and inelastic scattering (G4HadronHElasticPhysics, G4HadronPhysicsINCLXX, and G4IonElasticPhysics). The

59 simulations involved collisions of $^{28}$Si with an energy of 50 MeV/u and $^{12}$C particles in vacuum. The detector system consisted

60 of four supertelescope detectors. The simulation generated a dataset with more than four million particles. Fig. 2 shows the

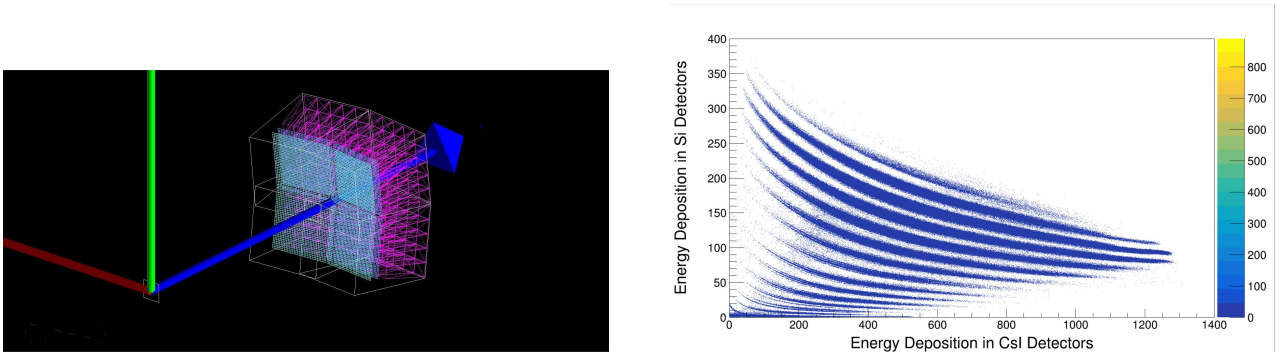61 structure of the detector system and the $\Delta$E-E two-dimensional histogram.



Fig. 2. The structure of the super telescope detector used in the Geant4 simulation and the $\Delta$E-E two-dimensional histogram from Geant4 simulation.

62 This study covers serval common machine learning algorithms, such as Support Vector Machines (SVM), Logistic Regression

63 (LR), and Bayesian classifiers. Ensemble learning algorithms based on tree structures and TabNet, a deep-learning algorithm,

64 were also used. The algorithms used in this study are briefly described below.

65 (a) MLP

66 Multi-layer perceptron (MLP) is a feed-forward neural network composed of multiple neurons, which is the basis and proto-

67 type of many artificial and deep learning neural networks.

68 (b) Random Forest

Random forest is an early tree-based ensemble learning algorithm with multiple decision trees [65]. This offers advantages of both decision trees and ensemble learning. Strong robustness and predictive ability are also advantageous.

(c) XgBoost

XgBoost is a tree-based ensemble learning algorithm proposed in 2016 [66] that is widely used in data mining, natural language processing, image recognition, and other fields. In general, XgBoost is a machine-learning algorithm with high efficiency, accuracy, flexibility, explainability, and scalability.

(d) LightGBM

LightGBM, a tree-based gradient boosting framework for ensemble learning, has been widely used in various applications [67]. Built on the gradient-boosting decision tree (GBDT) algorithm, LightGBM incorporates advanced techniques such as gradient-based one-sided sampling (GOSS) and histogram-based acceleration. These optimizations enabled faster training and lower memory consumption, making LightGBM an efficient and practical choice for machine learning tasks.

(e) CatBoost

CatBoost is a tree-based ensemble-learning algorithm developed by Yandex [68]. In terms of building a decision tree, compared with XgBoost and LightGBM, CatBoost can automatically process the category features of the data and automatically process the scaling of the data features without additional data processing. CatBoost adopts the same gradient-based splitting and feature selection strategies based on a greedy algorithm as XgBoost. CatBoost also automatically handles missing values in the data without additional data padding and has a certain robustness to noise and outliers.

Boosting-based ensemble learning algorithms such as XgBoost, LightGBM, and CatBoost are widely used in various fields. The basic process of these algorithms involves training multiple weak learners, assigning weights to training samples, and iteratively adjusting these weights based on the learner's performance. This iterative process aims to create a powerful ensemble model that is capable of accurate classification. Fig. 3 depicts the underlying structure of ensemble learning algorithms that employ the boosting method.

(f) TabNet

TabNet, which was introduced by Google in 2019, is a neural network structure explicitly designed for classification, prediction, and regression tasks involving tabular data [69]. Unlike traditional decision-tree-based machine learning algorithms, TabNet minimizes the need for preprocessing input data and can automatically learn the interdependencies among input features. It incorporates an attention transformer that uses an attention mechanism to select relevant feature vectors dynamically. Since its inception, TabNet has been widely adopted in various applications involving tabular data [70, 71].

Fig. 4 illustrates the procedure for applying the intelligent algorithms in this study. Training a classification model typically involves several steps.

(a) Data acquisition: Obtaining a dataset containing information about particle charge and mass, which can be from experimental or simulated data.

(b) Data preprocessing: Ensuring the quality and consistency of the data through noise removal, addressing missing data, and normalizing features.

(c) Data splitting: Dividing the dataset into training and testing sets. A training set was used to train the model and a test set was employed to evaluate the trained model. Random and stratified sampling are the commonly used methods.

(d) Feature engineering: Raw data is transformed, extracted, and selected to create informative and expressive feature sets.

(e) Algorithm selection: Suitable algorithms are chosen based on specific task requirements. The main task of the algorithm is mult-classification.
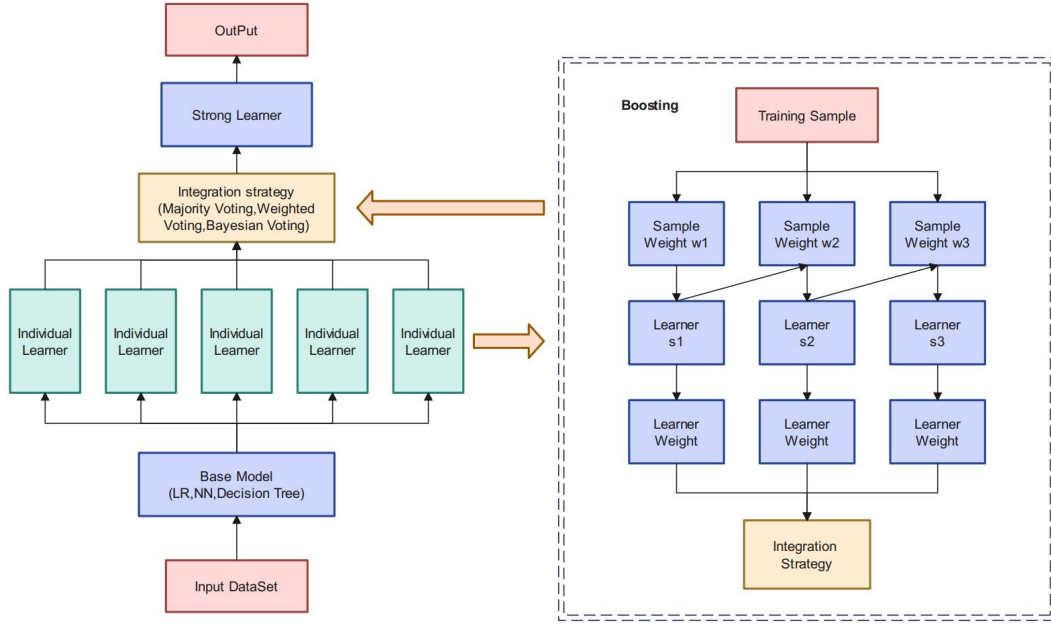
Fig. 3. The structure of ensemble learning and boosting method.

(f) Training and tuning of parameters: The algorithm's parameters can be tuned to improve the performance of the model further. Generally, each algorithm has a unique set of parameters that can be adjusted.

(g) The performance of the trained model is assessed using the testing dataset. Appropriate evaluation metrics were selected to evaluate the performance of the model for particle identification.

Based on the structure of NIMROD–ISiS, the dataset was initially split based on the ring number determined by the forward angle of the detector. Subsequently, the data were divided into two categories: telescope and supertelescope detectors. The Geant4 dataset was used for training and testing with machine learning and deep learning algorithms. After identifying the optimal algorithms, a subset of the detector data was used to evaluate the generalization ability of the algorithms.

Two classification strategies are adopted in this study.

(a) Using the algorithm to train and test charge and mass numbers, respectively.

(b) In classifying particle mass numbers, the particle's charge number was included as a part of the data features. From a logical perspective, this strategy is similar to traditional particle identification methods.

In a practical study, the experimental data exhibited a highly unbalanced distribution. Randomly extracting data can lead to disparate data category distributions among the training, validation, and test sets. This can lead to critical code errors and poor performance. Therefore, to address this problem, stratified sampling was employed as an alternative to random sampling.

## III.   RESULTS AND DISCUSSION

As the core task of particle identification involves multiple classifications, the use of suitable evaluation metrics for multiple classification algorithms is crucial. Common evaluation metrics include the accuracy, recall, precision, and f1-score [72–76]. These metrics help assess the performance of the algorithm from different aspects. The results of the classification task can be categorized into the following four types:
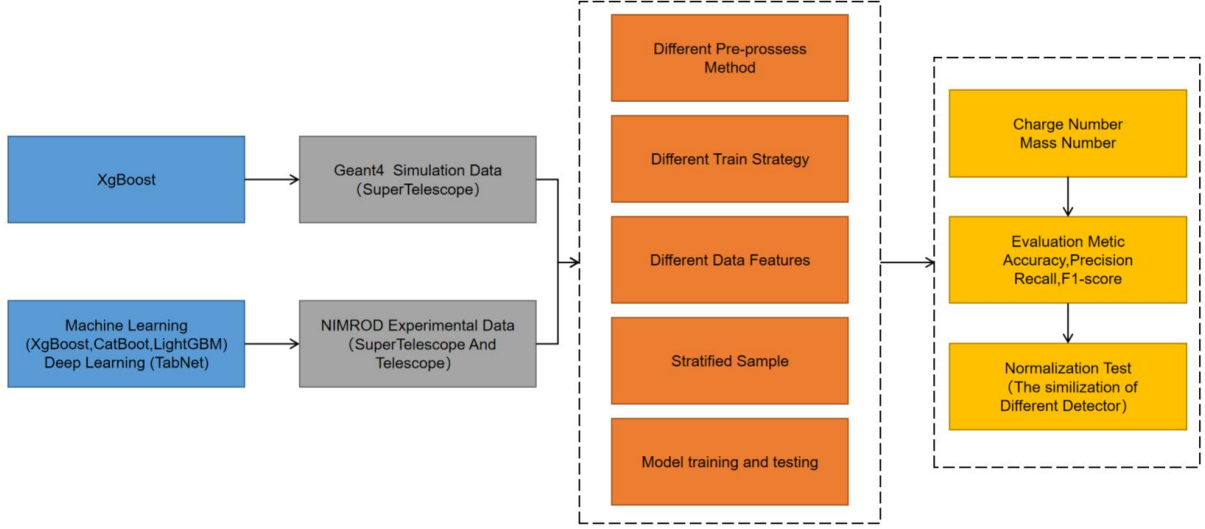
Fig. 4. The procedure of applying intelligent algorithms in this paper. The process is divided into three parts: datasets and algorithms, model training and testing, and evaluation of test results. Model training and testing is the most important part.

128 (a) Predict positive samples as positive. (TP)

129 (b) Predict negative samples as negative. (TN)

130 (c) Predict negative samples as positive. (FP)

131 (d) Predict positive samples as negative. (FN)

132 When evaluating the algorithm, the corresponding evaluation metrics were calculated using the classification results. The

133 equations are as follows:

$$Accuracy = (TP + TN)/(TP + TN + FP + FN) \tag{1}$$

$$Precision = (TP)/(TP + FP) \tag{2}$$

$$Recall = (TP)/(TP + FN) \tag{3}$$

$$F1 - score = 2 * Precision * Recall/(Precision + Recall) \tag{4}$$

141 Accuracy is defined as a measure of correctness. The precision measures the accuracy of a model in predicting positive

142 examples. The recall represents the coverage of positive samples that are correctly predicted. The f1-score is a compound

143 evaluation metric consisting of precision and recall.

144 Because the problem of positive and negative samples is extended to multiple categories in multi-classification tasks, meth-

145 ods for computing comprehensively evaluated metrics are essential. Commonly used strategies are the micro average, macro

146 average, and weighted average.

147 Macro-averaging calculates the average precision and recall of each class.

148 The micro-average ignores category differences and calculates the overall TP, FP, TN, and FN.

149 The weighted average is similar to the macro average, but uses category proportions as weights to calculate performance

150 metrics.

151 In particle identification, all generated particles have equal significance. Therefore, the macro average was chosen as the

152 calculation method for the evaluation metrics. The macro-average provides a balanced assessment across all classes and facili-

153 tates a comprehensive understanding of model performance. Because the mass and charge determined the particle category, the

154 charge and mass numbers were merged into a binary data format to calculate the evaluation metric.

155 The particles detected by the NIMROD–ISiS detector array can be categorized as light ions (with proton numbers ranging

156 from one to four) or heavy ions. Most heavy ions cannot penetrate the Si detector, whereas most light particles pass through

157 it. The disparity in the production yield between light particles and heavy ions during the reaction process leads to imbalances

158 in data distribution. The dataset was split based on whether the particles hit the CsI detector. This approach solves the data

159 imbalance problem in particle analysis and improves the algorithm performance. The XgBoost ensemble-learning algorithm

160 was selected for testing. The input data features were the total energy, energy deposition in the Si and CsI detectors, and the

161 detector position. The charge and mass numbers of the particles were used as data labels. Fig. 5 shows the results of the model

162 based on the telescope data. Table 1 shows the results of the model on the super-telescope data with two different classification
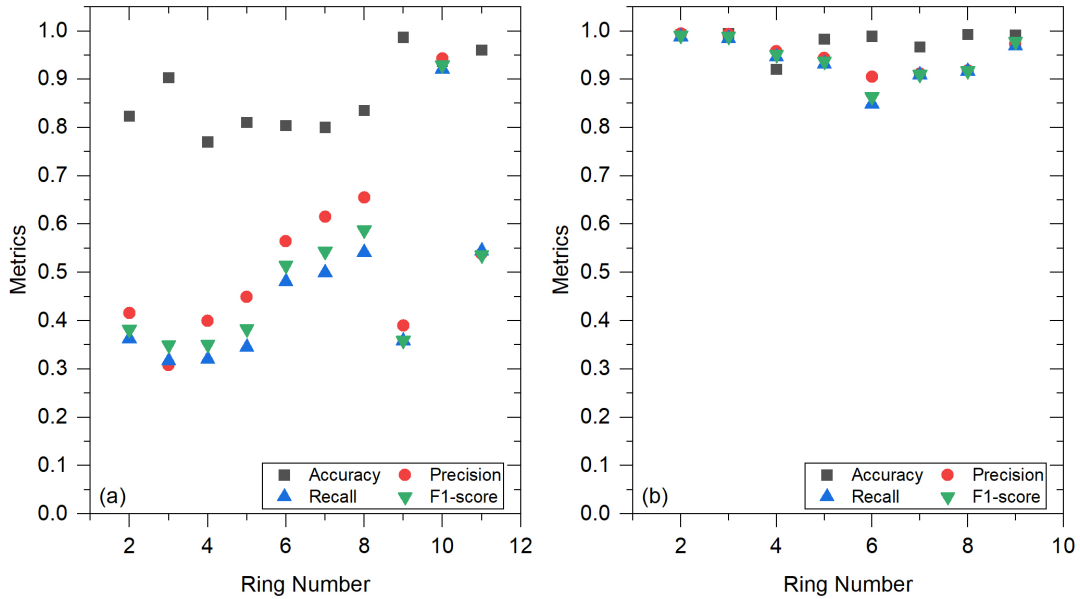
163 strategies.



Fig. 5. Test results of XgBoost on NIMROD–ISiS telescope data. Figures (a) and (b) show the results of XgBoost on particles w/o registrations on CsI detector. The latter results are better than the former.

164 The model performed well when tested on particles that were not registered on a CsI detector. The evaluation metrics for

165 each ring generally exceeded 0.85. The model performs better in identifying charge numbers than mass numbers. It also

166 achieves high accuracy for particles registered on CsI detector. However, their precision, recall, and f1-score were low. This

167 discrepancy is attributed to extreme data imbalance. Fig. 6 shows the mass distribution of ring 2. The mass distribution of

168 the particles registered on CsI detector was highly non-uniform. There was a significant difference between the categories with

169 the highest and lowest counts. The precise identification of rare categories is challenging for this model. As the evaluation

170 strategy uses a macro average, the evaluation metrics of the classifier are calculated as average values across all categories.

TABLE 1. Test results on NIMROD–ISiS SuperTelescope data.

| CsIE | Accuracy | Precision | Recall | F1-score | Label | Strategy |
|------|----------|-----------|--------|----------|-------|----------|
| Zero | 0.996 | 0.996 | 0.957 | 0.969 | Z | Independence |
| Zero | 0.934 | 0.874 | 0.844 | 0.856 | A | Independence |
| Zero | 0.932 | 0.893 | 0.877 | 0.883 | Z+A | Independence |
| Zero | 0.997 | 0.997 | 0.958 | 0.969 | Z | FirstZ,SecondA |
| Zero | 0.966 | 0.908 | 0.893 | 0.9 | A | FirstZ,SecondA |
| Zero | 0.964 | 0.924 | 0.916 | 0.919 | Z+A | FirstZ,SecondA |
| Non-Zero | 0.974 | 0.473 | 0.402 | 0.425 | Z | Independence |
| Non-Zero | 0.892 | 0.3 | 0.244 | 0.261 | A | Independence |
| Non-Zero | 0.87 | 0.316 | 0.247 | 0.266 | Z+A | Independence |
| Non-Zero | 0.974 | 0.468 | 0.401 | 0.425 | Z | FirstZ,SecondA |
| Non-Zero | 0.903 | 0.326 | 0.274 | 0.289 | A | FirstZ,SecondA |
| Non-Zero | 0.881 | 0.335 | 0.272 | 0.295 | Z+A | FirstZ,SecondA |

Thus, the performance of categories with small percentages significantly affected the overall evaluation metrics. The training strategies for charge and mass numbers did not show any significant differences. Including the charge number as an additional data feature did not effectively improve the identification ability of the model for particles in deficient quantities. If the model fails to precisely predict the charge number of the particles, the accuracy of the mass number identification is also affected.
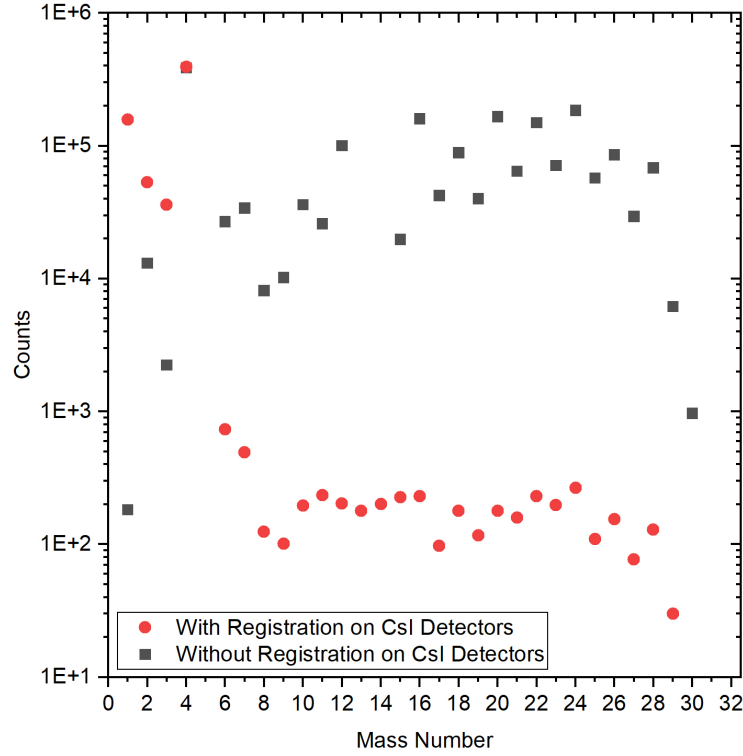


Fig. 6. Mass number distribution for events w/o registration on CsI detectors in ring 2. The mass number distribution of particles without registration on CsI detector is well balanced, with sample sizes exceeding 1000 in most categories. Among particles with registration on CsI detector, most heavy ions count around 100 occurrences.

To address this problem, the following methods have been proposed:

(a) Algorithm parameter optimization: Refining algorithm parameters (reducing the learning rate, increasing iteration numbers, expanding tree depth, etc.) to improve accuracy, precision, and recall. However, adjusting the parameters alone had a limited

178 impact on the categories with limited samples, even when distinct weights were assigned to each category.

179 (b) Data category adjustment: The imbalance ratio can be reduced by eliminating data categories that comprise only a few or a
180 few dozen samples.

181 (c) Exploration of data pre-processing methods: Trying out different approaches, including normalization, standardization, or
182 no data pre-processing.

183 The most effective solution to the severe shortage of samples in specific categories is to include additional data. This reduces
184 the imbalance ratio and thus improves the accuracy. For instance, in ring 10, each category had over 20,000 samples and the
185 imbalance ratio was only 5:1. XgBoost performed excellently, with the evaluation metrics for each type exceeding 0.9.

186 Other factors, such as detector position and hardware conditions, such as temperature and electronic signal drift, can cause
187 scaling issues, thus affecting algorithm accuracy. To address this issue, Geant4 was used to simulate the experiment and detector
188 performance, enabling a focused research to address the imbalance issue.

189 In Geant4, the total particle energy, time of flight (ToF), kinetic energy before entering the detector, detector position, and
190 particle deposition energy (Eabs) were selected as input data features. Testing with XgBoost demonstrated that the additional
191 data features alleviated the data imbalance problem, resulting in excellent performance.

192 To confirm that this is not limited to XgBoost alone, a comparison test was conducted with other machine-learning and deep-
193 learning algorithms. The test results ( Fig. 7) confirmed the earlier findings. Tree-based machine learning algorithms, such
194 as XgBoost and deep learning TabNet demonstrated excellent performance, whereas traditional machine learning algorithms,
195 such as LR, SVM, and Bayesian classifiers, exhibited poor performance.

196 These results validate the effectiveness of the proposed approach in mitigating data imbalances and highlight the superiority
197 of tree-based machine learning and deep learning algorithms in addressing this challenge.
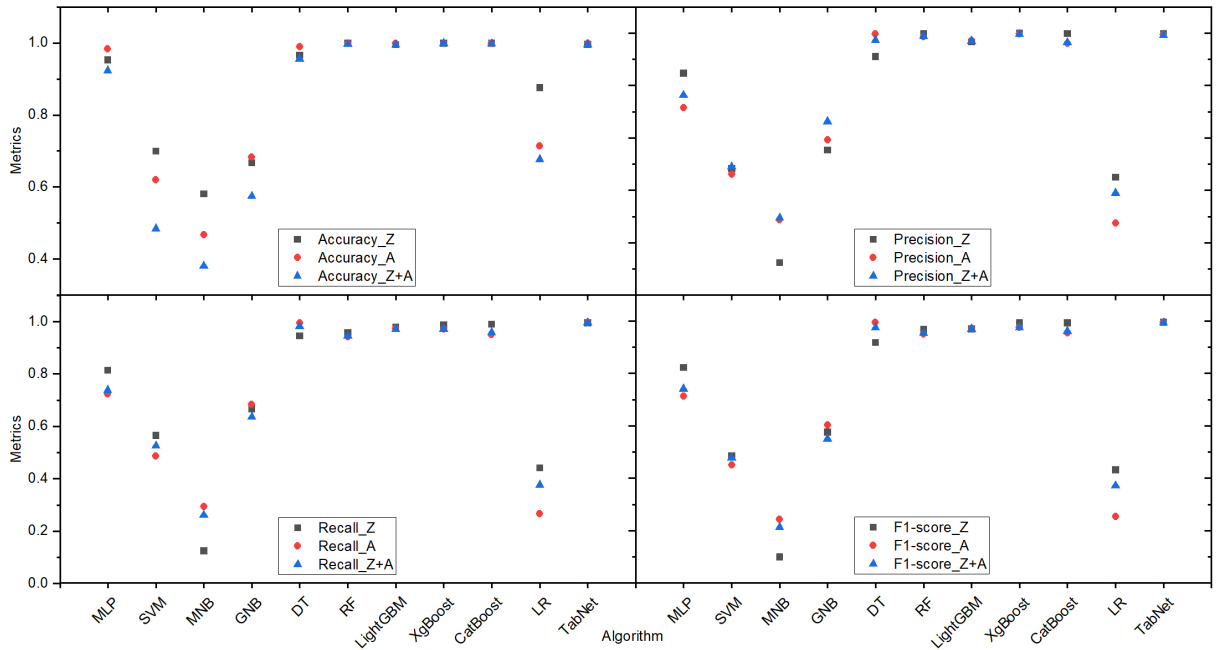


Fig. 7. Test results of machine learning and deep learning algorithms on the Geant4 dataset. SVM, MNB, GNB, and LR perform poorly. MLP has relatively low precision and recall. Ensemble learning algorithms such as XgBoost and deep learning algorithm TabNet perform well.

198 Subsequently, the algorithms were evaluated using only the energy deposition as the data feature. XgBoost, LightGBM,

CatBoost, and TabNet, which exhibited promising performances in prior tests, were selected for this assessment. The results demonstrated that each algorithm showed decreased accuracy in predicting the particle mass number.

Based on these observations, a series of additional data features were selected for comparative analysis. Numerous tests have shown that particle flight time is important for improving the accuracy of the algorithm. Although this feature alone proved insufficient for charge and mass number identification, its combination with the energy deposition feature enhanced the identification capabilities of the algorithm. Detailed descriptions of the corresponding test results are presented in Table 2 and Table 3.

TABLE 2. Classification results from Geant4 simulation data, with independent training on charge and mass numbers.

| Algorithm | Accuracy | Precision | Recall | F1-score | Feature |
|---|---|---|---|---|---|
| XgBoost | 0.127 | 0.045 | 0.079 | 0.05 | ToF |
| XgBoost | 0.862 | 0.863 | 0.827 | 0.839 | Eabs |
| LightGBM | 0.828 | 0.821 | 0.795 | 0.804 | Eabs |
| CatBoost | 0.836 | 0.804 | 0.765 | 0.771 | Eabs |
| TabNet | 0.813 | 0.837 | 0.762 | 0.791 | Eabs |
| XgBoost | 0.97 | 0.986 | 0.963 | 0.971 | Eabs,ToF |
| LightGBM | 0.947 | 0.95 | 0.936 | 0.943 | Eabs,ToF |
| CatBoost | 0.948 | 0.948 | 0.914 | 0.926 | Eabs,ToF |
| TabNet | 0.971 | 0.99 | 0.976 | 0.983 | Eabs,ToF |

TABLE 3. Classification results from Geant4 simulation data, where charge number is one of the data features of dataset defined by mass number.

| Algorithm | Accuracy | Precision | Recall | F1-score | Feature |
|---|---|---|---|---|---|
| XgBoost | 0.127 | 0.056 | 0.079 | 0.051 | ToF |
| XgBoost | 0.87 | 0.878 | 0.839 | 0.852 | Eabs |
| LightGBM | 0.85 | 0.848 | 0.812 | 0.82 | Eabs |
| CatBoost | 0.83 | 0.798 | 0.752 | 0.757 | Eabs |
| TabNet | 0.828 | 0.854 | 0.794 | 0.821 | Eabs |
| XgBoost | 0.971 | 0.987 | 0.965 | 0.972 | Eabs,ToF |
| LightGBM | 0.952 | 0.949 | 0.943 | 0.945 | Eabs,ToF |
| CatBoost | 0.948 | 0.947 | 0.906 | 0.918 | Eabs,ToF |
| TabNet | 0.952 | 0.985 | 0.948 | 0.963 | Eabs,ToF |

The final phase of the study involved a comprehensive investigation of the generalization ability of the algorithms. Unlike previous tests involving data from all detectors, this phase focuses on training models using a specific subset of detectors, reserving the remaining data for testing. The data features include the time-of-flight (ToF) and energy deposition. Various data preprocessing techniques, including normalization and standardization, were explored during the testing phase. Before model training, datasets from specific detectors were normalized and standardized using the MinMaxScaler and StandardScaler methods from the sklearn.preprocessing package in Python. These methods were also used to test data from other detectors before evaluating the trained model. However, these methods have a significantly negative impact on generalization ability. Therefore, no data preprocessing was performed. The results are shown in Fig. 8.

The performance of the algorithms was excellent. The evaluation metrics of TabNet and XgBoost are mostly over 0.9 for all detector data (Fig. 9). These findings establish the efficacy of training models with robust generalization abilities even in situations with limited data availability. Overall, these results highlight the advantages and effectiveness of machine learning and deep learning algorithms, and demonstrate their potential for practical applications.

Inspired by these findings, a similar study of data similarity was conducted on specific rings of NIMROD–ISiS. The input
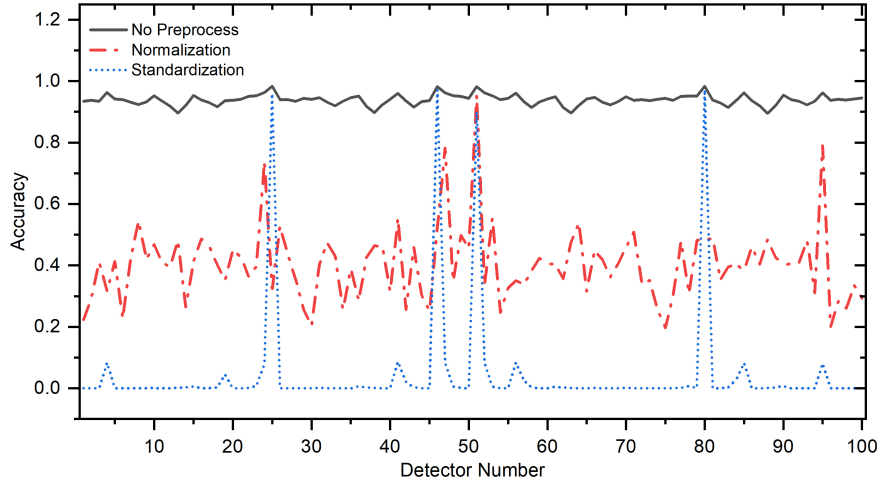
Fig. 8. Generalization ability test results of the XgBoost algorithm on the Geant4 dataset using different data preprocessing methods. It can be noticed that both normalization and standardization severely reduce the model's generalization ability.
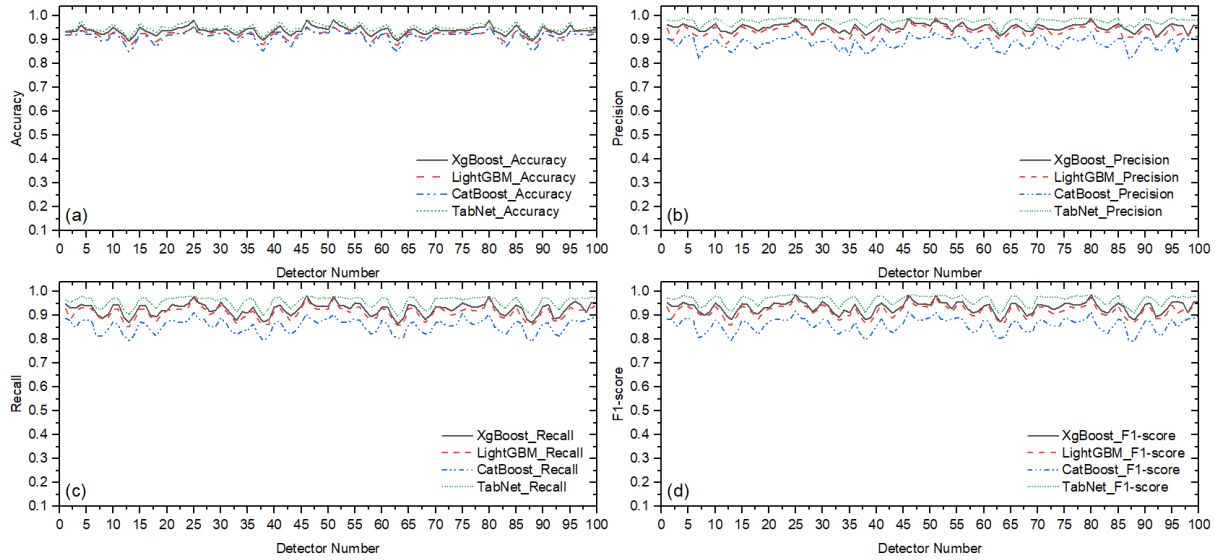


Fig. 9. The test results of generalizability ability test of XgBoost, CatBoost, LightGBM, and TabNet. Figures (a), (b), (c), and (d), respectively, show their accuracy, precision, recall, and f1-score curves. The evaluation metrics of all algorithms exceed 0.8. The evaluation metrics of XgBoost and TabNet are mostly over 0.9. TabNet shows better generalization ability than ensemble learning algorithms.

219 data features included total energy, energy deposition, and detector position. Through testing, it was discovered that, depending
220 on the similarity of the data, the detectors of NIMROD–ISiS can be divided into groups. Taking the data (particles registered
221 on CsI detector) from ring 9 as an example, ring 9 can be further divided into two groups of detectors (Fig. (a) and Fig. (b) of
222 Fig. 10). The results depicted in Fig. 10 provide valuable information on the patterns and characteristics of the NIMROD–ISiS
223 detector. Moreover, they contribute to the optimization of algorithms and the improvement of data-processing methods. These
224 findings also have significant implications for the study and enhancement of the detector array design and performance. Overall,
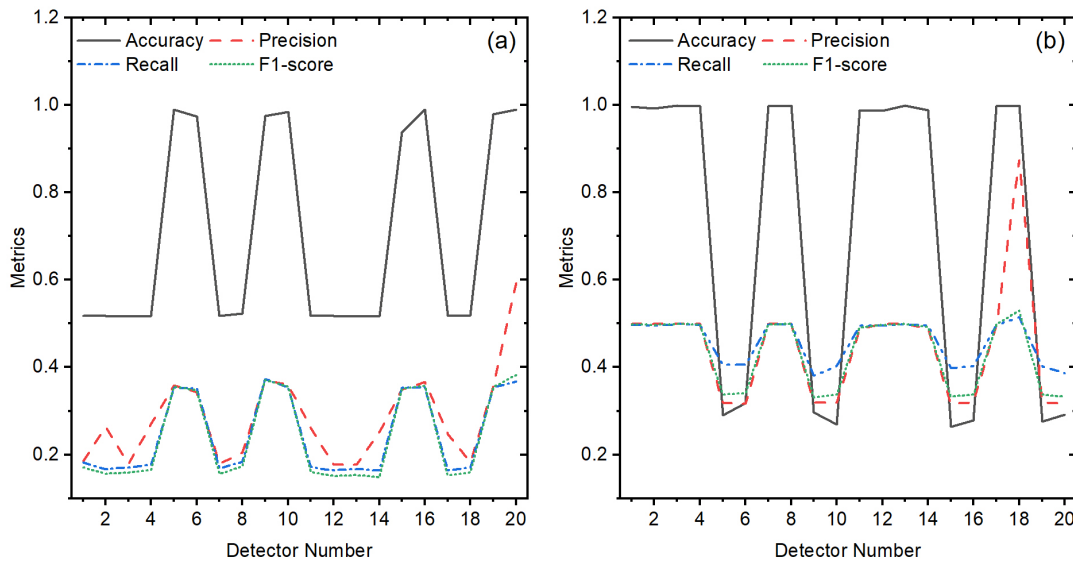225 the findings have practical implications.

Fig. 10. In the tests carried out on the NIMROD–ISiS ring 9, the predictions for the other detector data can be categorized into two scenarios, Fig. (a) and Fig. (b). A high degree of similarity can be observed between the detectors on ring 9.

## IV. CONCLUSION

Particle identification in machine learning (ML) is an integrated problem. Researchers must consider various factors including data selection, partitioning, feature engineering, preprocessing, algorithm selection, and parameter tuning. Traditional particle identification methods require significant manual effort and are limited by researchers' experience and available time. Our study aims to develop a universal and adaptable particle identification model that assists in manual processes. Although achieving 100% accuracy may not be possible, ensemble learning algorithms have meaningful results, especially XgBoost. The conclusions are as follows:

First, intelligent algorithms, particularly tree-based ensemble learning algorithms, can effectively identify particles in heavy-ion collisions at low and intermediate energies. This offers a viable alternative to traditional methods.

Secondly, addressing data imbalances is crucial for particle identification. Severe data imbalances significantly affected the results. The solutions include ensuring sufficient data for a balanced distribution, adding additional data features beyond particle energy deposition, and constructing different identification models based on the detector structure.

Third, training a specialized particle identification model using the existing data reduces the time and resources required for traditional particle identification. Laboratories conducting long-term, large-scale heavy-ion collision experiments can be beneficial. This paves the way for the development of a professional particle identification software.

Finally, machine-learning algorithms can be used to study detector similarity, particularly in large-scale detector arrays with complex structures.

Combinations of supervised and unsupervised learning approaches should be explored in future studies. Other physics software such as NpTool [77] will also be used to simulate the experiments. NpTool is known for its efficient project management and simulation of various sophisticated detector arrays.

Because Geant4 simulations are time consuming and resource intensive, there is a need to explore alternative approaches for generating particle collision data. Generative Adversarial Networks (GAN) [78] and Variational Autoencoders (VAE) [79] have shown promise in generating simulated data for detectors in the field of high-energy physics [80–83]. Utilizing GAN and VAE

can reduce the time and resources required for massive amounts of simulated data, thereby making the process more efficient and accessible.

Building on the excellent performance of TabNet, further investigations will include exploring additional deep-learning algorithms, such as DeepGBM [84] and GrowNet [85]. Moreover, we attempted to change the existing ensemble learning algorithm into a multi-output algorithm to classify the mass and charge numbers simultaneously. Our research aims to enhance the understanding of the detector system in sophisticated experiments, which can be used to explore interesting clustering phenomena in nuclei [86–89].

## V. ACKNOWLEDGEMENTS

[1] A. Kalweit, Particle identification in the ALICE experiment. J. Phys. G. Nucl. Partic. **38**, 124073 (2011).

[2] C. Zampolli, Particle identification with the ALICE detector at the LHC. (2012).

[3] P. Križan, Particle identification at Belle II. J. INSTRUM. **9**, C07018 (2014).

[4] Ł.K. Graczykowski, M. Jakubowska, K.R. Deja et al., Using machine learning for particle identification in ALICE. J. INSTRUM. **17**, C07016 (2022).

[5] P. Calafiura, S. Farrell, H. Gray et al., TrackML: a high energy physics particle tracking challenge. IEEE 14th International Conference on E-Science (e-Science) 344-344 (2018).

[6] C. Tüysüz, F. Carminati, B. Demirköz et al., Particle track reconstruction with quantum algorithms. Epj. Web. Conf. **245**, 09013 (2020).

[7] O. Bakina, D. Baranov, I. Denisenko et al., Deep learning for track recognition in pixel and strip-based particle detectors. J. INSTRUM. **17**, P12023 (2022).

[8] P. Goncharov, E. Schavelev, A. Nikolskaya et al., Ariadne: PyTorch library for particle track reconstruction using deep learning. AIP Conf. Proc. **2377**, (2021).

[9] T.Q. Chen, T. He, Higgs boson discovery with boosted trees. JMLR Work. Conf. Proc. **42**, 69-80 (2015).

[10] M. Azhari, A. Abarda, B. Ettaki et al., Higgs boson discovery using machine learning methods with PySpark. Procedia Comput. Sci. **170**, 1141-1146 (2020).

[11] C. Adam-Bourdarios, G. Cowan, C. Germain-Renaud et al., The Higgs machine learning challenge. J. Phys.: Conf. Ser. **664**, 072015 (2015).

[12] S.R. Ahmad, Technical report of participation in Higgs boson machine learning challenge. (2015).

[13] A.E. Phoboo, Machine learning wins the Higgs challenge. (CERN Bulletin, 2014), https://cds.cern.ch/journal/CERNBulletin/2014/49/News%20Articles/1972036. Accessed November 16, 2023.

[14] University of California, Irvine. ML Physics Portal, http://mlphysics.ics.uci.edu/. Accessed November 16, 2023.

[15] M.J. Fenton, A. Shmakov, TW. Ho et al., Permutationless many-jet event reconstruction with symmetry preserving attention networks. Phys. Rev. D **105**, 112008 (2022).

[16] C. Shimmin, P. Sadowski, P. Baldi et al., Decorrelated jet substructure tagging using adversarial neural networks. Phys. Rev. D **96**, 074034 (2017).

[17] P. Baldi, P. Sadowski, D. Whiteson, Searching for exotic particles in high-energy physics with deep learning. Nat. Commun. **5**, 4308 (2014).

[18] W.B. He, J.J. He, R. Wang et al., Machine learning applications in nuclear physics (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252004 (2022).

[19] W.B. He, Q.F. Li, Y.G. Ma et al., Machine learning in nuclear physics at low and intermediate energies. Sci. China Phys. Mech. Astron. **66**, 282001 (2023).

[20] M. Zhou, Y.Q. Luo, H.C. Song, Applications of machine learning in relativistic heavy ion physics (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252002 (2022).

[21] W.B. He, Y.G. Ma, L.G. Pang et al., High-energy nuclear physics meets machine learning. Nucl. Sci. Tech. **34**, 88 (2023).

[22] T.L. Zhao, H.F. Zhang, Neural network approach to improve the quality of atomic nuclei (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252008 (2022).

[23] M.R. Mumpower, T.M. Sprouse, A.E. Lovell et al., Physically interpretable machine learning for nuclear masses. Phys. Rev. C **106**, L021301 (2022).

[24] Z.M. Niu, H.Z. Liang, Nuclear mass predictions with machine learning reaching the accuracy required by $r$-process studies. Phys. Rev. C **106**, L021303 (2022).

[25] L. Neufcourt, Y.C. Cao, W. Nazarewicz et al., Bayesian approach to model-based extrapolation of nuclear observables. Phys. Rev. C **98**, 034318 (2018).

[26] Z.P. Gao, Y.J. Wang, H.L. Lü et al., Machine learning the nuclear mass. Nucl. Sci. Tech. **32**, 109 (2021).

[27] J.Q. Ma, Z.H. Zhang, Improved phenomenological nuclear charge radius formulae with kernel ridge regression. Chinese Phys. C **46**, 074105 (2022).

[28] X.X. Dong, R. An, J.X. Lu et al., Nuclear charge radii in Bayesian neural networks revisited. Phys. Lett. B **838**, 137726 (2023).

[29] R. Utama, W.C. Chen, J. Piekarewicz, Nuclear charge radii: density functional theory meets Bayesian neural networks. J. Phys. G. Nucl. Partic. **43**, 114002 (2016).

[30] S.J. Tao, L.F. Zhang, Q.Y. Zhang et al., Improved naive Bayesian probability classifier in nuclear charge radius prediction (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252009(2022).

[31] Y.F. Ma, C. Su, J. Liu et al., Predictions of nuclear charge radii and physical interpretations based on the naive Bayesian probability classifier. Phys. Rev. C **101**, 014304 (2020).

[32] Z.M. Niu, H.Z. Liang, B.H. Sun et al., Predictions of nuclear $\beta$-decay half-lives with machine learning and their impact on $\gamma$-process nucleosynthesis. Phys. Rev. C **99**, 064307 (2019).

[33] J.M. Munoz, S. Akkoyun, Z.P. Reyes et al., Predicting $\beta$-decay energy with machine learning. Phys. Rev. C **107**, 034308 (2023).

[34] Z.Y. Yuan, D. Bai, Z.Z. Ren et al., Theoretical predictions on $\alpha$-decay properties of some unknown neutron-deficient actinide nuclei using machine learning. Chinese Phys. C **46**, 024101 (2022).

[35] X.D. Bu, D. Wu, C.L. Bai, Prediction of $\alpha$-decay half-lives for superheavy nuclei based on neural network (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252005 (2022).

[36] P. Li, J.H. Bai, Z.M. Niu et al., $\beta$-decay half-lives studied using neural network method (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252006 (2022).

[37] N.J. Costiris, E. Mavrommatis, K.A. Gernoth et al., Decoding $\beta$-decay systematics: a global statistical model for $\beta$- half-lives. Phys. Rev. C **80**, 044332 (2009).

[38] R. Wang, Y.G. Ma, R. Wada et al., Nuclear liquid-gas phase transition with machine learning. Phys. Research **2**, 043202 (2020).

[39] D. Peng, H.L. Wei, J. Pu et al., Bayesian neural network prediction methods for fragment cross sections in proton-induced spallation reactions (in Chinese). Sci. Sin.-Phys. Mech. Astron. **52**, 252012 (2022).

[40] B.C. Wang, M.T. Qiu, W. Chen et al., Machine learning-based analyses for total ionizing dose effects in bipolar junction transistors. Nucl. Sci. Tech. **33**, 131 (2022).

[41] Y.B. Yu, G.F. Liu, W. Xu et al., Research on tune feedback of the Hefei Light Source II based on machine learning. Nucl. Sci. Tech. **33**, 28 (2022).

[42] Y.D. Song, R. Wang, Y.G. Ma et al., Determining the temperature in heavy-ion collisions with multiplicity distribution. Phys. Lett. B **814**, 136084 (2021).

[43] Q.F. Song, L. Zhu, J. Su, Target dependence of isotopic cross sections in the spallation reactions $^{238}$U $+ p$ , $d$ and $^9$Be at 1 $A$GeV. Chinese Phys. C **46**, 074108 (2022).

[44] H.K. Wu, Y.J. Wang, Y.M. Wang et al., Machine learning method for $^{12}$ C event classification and reconstruction in the active target time-projection chamber. Nucl. Instrum. Meth. A (2023).

[45] F.P. Li, Y.J. Wang, Z.P. Gao et al., Application of machine learning in the determination of impact parameter in the $^{132}$Sn $+$ $^{124}$Sn system. Phys. Rev. C **104**, 034608 (2021).

[46] Z.Y. Li, Z. Qian, J.H. He et al., Improvement of machine learning-based vertex reconstruction for large liquid scintillator detectors with multiple types of PMTs. Nucl. Sci. Tech. **33**, 93 (2022).

[47] H. Arahmane, EM. Hamzaoui, Y.B. Maissa et al., Neutron-gamma discrimination method based on blind source separation and machine learning. Nucl. Sci. Tech. **32**, 18 (2021).

[48] J. Collado, J.N. Howard, T. Faucett et al., Learning to identify electrons. Phys. Rev. D **103**, 116028 (2021).

[49] L. de Oliveira, B. Nachman, M. Paganini, Electromagnetic showers beyond shower shapes. Nucl. Instrum. Meth. A **951**, 162879 (2020).

[50] P. Baldi, K. Bauer, C. Eng et al., Jet substructure classification in high-energy physics with deep neural networks. Phys. Rev. D **93**, 094034 (2016).

[51] C. Fanelli, J. Pomponi, DeepRICH: learning deeply Cherenkov detectors. Mach. Learn.: Sci. Technol. **1**, 015010 (2020).

[52] E. Cisbani, A. Del Dotto, C. Fanelli et al., AI-optimized detector design for the future Electron-Ion Collider: the dual-radiator RICH case. J. INSTRUM. **15**, P05009 (2020).

[53] S. Carboni, S. Barlini, L. Bardelli et al., Particle identification using the $\Delta$E–E technique and pulse shape discrimination with the silicon detectors of the FAZIA project. Nucl. Instrum. Meth. A **664**, 251-263 (2012).

[54] W. Klempt, Review of particle identification by time of flight techniques. Nucl. Instrum. Meth. A **433**, 542-553 (1999).

[55] Y.G. Ma, Effects of $\alpha$-clustering structure on nuclear reaction and relativistic heavy-ion collisions. Nuclear Techniques **46**, 080001 (2023).

[56] J.J. He, W.B. He, Y.G. Ma et al., Machine-learning-based identification for initial clustering structure in relativistic heavy-ion collisions. Phys. Rev. C **104**, 044902 (2021).

[57] Y.G. Ma, S. Zhang, Influence of nuclear structure in relativistic heavy-ion collisions. Handbook of Nuclear Physics 1-30 (2022).

[58] X.G. Cao, E.J. Kim, K. Schmidt et al., Examination of evidence for resonances at high excitation energy in the 7 $\alpha$ disassembly of $^{28}$Si. Phys. Rev. C **99**, 014606 (2019).

[59] X.G. Cao, E.J. Kim, K. Schmidt et al., $\alpha$ and $\alpha$ conjugate fragment decay from the disassembly of $^{28}$Si at very high excitation energy. JPS Conf. Proc. 010038 (2020).

[60] X.G. Cao, E.J. Kim, K. Schmidt et al., Evidence for resonances in the 7 $\alpha$ disassembly of $^{28}$Si. AIP Conf. Proc. **2038**, 020021 (2018).

[61] P. Adamson, M. Youngs, Machine learning: potential application for particle identification. 2019 Fall Meeting of the APS Division of Nuclear Physics (2019).

[62] S. Wuenschel, K. Hagel, R. Wada et al., NIMROD–ISiS, a versatile tool for studying the isotopic degree of freedom in heavy ion collisions. Nucl. Instrum. Meth. A **604**, 578-583 (2009).

[63] R. Wada, S. Wuenschel, K. Hagel et al., A $4\pi$ detector array, NIMROD-ISIS. Nucl. Phys. News **24**, 28-33 (2014).

[64] S. Agostinelli, J. Allison, K. Amako et al., GEANT4—a simulation toolkit. Nucl. Instrum. Meth. A **506**, 250-303 (2003).

[65] L. Breiman, Random forests. Mach. Learn. **45**, 5-32 (2001).

[66] T.Q. Chen, C. Guestrin, XGBoost: a scalable tree boosting system. Proceedings of the 22nd ACM Sigkdd International Conference on Knowledge Discovery and Data Mining 785-794 (2016).

[67] G.L. Ke, Q. Meng, T. Finley et al., LightGBM: a highly efficient gradient boosting decision tree. Adv. Neur. In. **30**, (2017).

[68] A.V. Dorogush, V. Ershov, A. Gulin, CatBoost: gradient boosting with categorical features support. (2018).

[69] S.Ö. Arik, T. Pfister, Tabnet: attentive interpretable tabular learning. In Proceedings of the AAAI Conference on Artificial Intelligence **35**, 6679-6687 (2021).

[70] J.Z. Yan, T.Y. Xu, Y.C. Yu et al., Rainfall forecast model based on the tabnet model. Water **13**, 1272 (2021).

[71] R. Asencios, C. Asencios, E. Ramos, Profit scoring for credit unions using the multilayer perceptron, XGBoost and TabNet algorithms: evidence from Peru. Expert Syst. Appl. **213**, 119201 (2023).

[72] B. Juba, H.S. Le, Precision-recall versus accuracy and the role of large data sets. Proceedings of the AAAI Conference on Artificial Intelligence **33**, 4039-4048 (2019).

[73] N. Japkowicz, Assessment metrics for imbalanced learning. Imbalanced Learning: Foundations, Algorithms, and Applications 187-206 (2013).

[74] M. Grandini, E. Bagli, G. Visani, Metrics for multi-class classification: an overview. (2020).

[75] M. Hossin, M.N. Sulaiman, A review on evaluation metrics for data classification evaluations. IJDKP. **5**, 1 (2015).

[76] E. Mortaz, Imbalance accuracy metric for model selection in multi-class imbalance classification problems. Knowl.-Based Syst. **210**, 106490 (2020).

[77] A, Matta, P. Morfouace, N. de Séréville et al., NPTool: a simulation and analysis framework for low-energy nuclear physics experiments. J. Phys. G. Nucl. Partic. **43**, 045113 (2016).

[78] I. Goodfellow, J. Pouget-Abadie, M. Mirza et al., Generative adversarial networks. Commun. ACM **63**, 139-144 (2020).

[79] D.P. Kingma, M. Welling, Auto-encoding variational bayes. (2013).

[80] D. Derkach, N. Kazeev, F. Ratnikov et al., Cherenkov detectors fast simulation using neural networks. Nucl. Instrum. Meth. A **952**, 161804 (2020).

[81] M. Paganini, L. de Oliveira, B. Nachman, Accelerating science with generative adversarial networks: an application to 3D particle showers in multilayer calorimeters. Phys. Rev. Lett. **120**, 042003 (2018).

[82] D. Salamani, S. Gadatsch, T. Golling et al., Deep generative models for fast shower simulation in ATLAS. IEEE 14th International Conference on E-Science (e-Science) 348 (2018).

[83] D. Belayneh, F. Carminati, A. Farbin et al., Calorimetry with deep learning: particle simulation and reconstruction for collider physics. Eur. Phys. J. C **80**, 1-31 (2020).

[84] G.L. Ke, Z.H. Xu, J. Zhang et al., DeepGBM: a deep learning framework distilled by GBDT for online prediction tasks. Proceedings of the 22nd ACM Sigkdd International Conference on Knowledge Discovery and Data Mining 384-394 (2019).

[85] S. Badirli, X.Q. Liu, Z.M. Xing et al., Gradient boosting neural networks: Grownet. (2020).

[86] W.B. He, X.G. Cao, Y.G. Ma et al., Application of EQMD model to researches of nuclear exotic structures. Nuclear Techniques **37** (2014).

[87] X.G. Cao, Y.G. Ma, Progress of theoretical and experimental studies on $\alpha$ cluster structures in light nuclei. Chinese. Sci. Bull. **60**, 1557-1564 (2015).

[88] W.B. He, Y.G. Ma, X.G. Cao et al., Dipole oscillation modes in light $\alpha$-clustering nuclei. Phys. Rev. C **94**, 014301 (2016).

[89] W.B. He, Y.G. Ma, X.G. Cao et al., Giant dipole resonance as a fingerprint of $\alpha$ clustering configurations in $^{12}$C and $^{16}$O. Phys. Rev. Lett. **113**, 032506 (2014).